

An Automatic Finite-Sample Robustness Metric: Can Dropping a Little Data Make a Big Difference?

Ryan Giordano (rgiordan@mit.edu)¹
January 2022

¹With coauthors Rachael Meager (LSE) and Tamara Broderick (MIT)

Dropping data: Motivation

More data & cheaper computation \Rightarrow

Statistical analyses are playing larger roles in decision making.

Decisions are important: We want **trustworthy** conclusions.

Data / models not always perfect: We want **robust** conclusions.

Would you be concerned if you could **reverse your conclusion** by removing a **small proportion** (say, 0.1%) of your data?

Running example: Angelucci et al. [2015], a randomized controlled trial study of the efficacy of microcredit based on 16,560 data points.

We can reverse the studies qualitative conclusions by removing 15 observations ($< 0.1\%$ of the data).

How do we find sets of influential points? Difficult in general!

We provide a **automatic approximation** with finite-sample guarantees.

Studying the approximation reveals the causes of non-robustness.

Dropping data: Mexico Microcredit

Consider Angelucci et al. [2015], a randomized controlled trial study of the efficacy of microcredit in Mexico based on 16,560 data points.

The variable “Beta” estimates the effect of microcredit in US dollars.

	Beta (SE)
Original result	-4.55 (5.88)

Original conclusion:

There is no evidence that microcredit is effective.

Dropping data: Mexico Microcredit

Consider Angelucci et al. [2015], a randomized controlled trial study of the efficacy of microcredit in Mexico based on 16,560 data points.

The variable “Beta” estimates the effect of microcredit in US dollars.

	Left out points	Beta (SE)
Original result	0	-4.55 (5.88)
Change sign	1	0.4 (3.19)

Original conclusion:

There is no evidence that microcredit is effective.

Dropping data: Mexico Microcredit

Consider Angelucci et al. [2015], a randomized controlled trial study of the efficacy of microcredit in Mexico based on 16,560 data points.

The variable “Beta” estimates the effect of microcredit in US dollars.

	Left out points	Beta (SE)
Original result	0	-4.55 (5.88)
Change sign	1	0.4 (3.19)
Change significance	14	-10.96 (5.57)

Original conclusion:

There is no evidence that microcredit is effective.

Dropping data: Mexico Microcredit

Consider Angelucci et al. [2015], a randomized controlled trial study of the efficacy of microcredit in Mexico based on 16,560 data points.

The variable “Beta” estimates the effect of microcredit in US dollars.

	Left out points	Beta (SE)
Original result	0	-4.55 (5.88)
Change sign	1	0.4 (3.19)
Change significance	14	-10.96 (5.57)
Change sign and significance	15	7.03 (2.55)

Original conclusion:

There is no evidence that microcredit is effective.

Dropping data: Mexico Microcredit

Consider Angelucci et al. [2015], a randomized controlled trial study of the efficacy of microcredit in Mexico based on 16,560 data points.

The variable “Beta” estimates the effect of microcredit in US dollars.

	Left out points	Beta (SE)
Original result	0	-4.55 (5.88)
Change sign	1	0.4 (3.19)
Change significance	14	-10.96 (5.57)
Change sign and significance	15	7.03 (2.55)

Original conclusion:

There is no evidence that microcredit is effective.

Potential conclusions after data dropping:

The effect of microcredit is positive (negative) & statistically significant.

Dropping data: Mexico Microcredit

Consider Angelucci et al. [2015], a randomized controlled trial study of the efficacy of microcredit in Mexico based on 16,560 data points.

The variable “Beta” estimates the effect of microcredit in US dollars.

	Left out points	Beta (SE)
Original result	0	-4.55 (5.88)
Change sign	1	0.4 (3.19)
Change significance	14	-10.96 (5.57)
Change sign and significance	15	7.03 (2.55)

Original conclusion:

There is no evidence that microcredit is effective.

Potential conclusions after data dropping:

The effect of microcredit is positive (negative) & statistically significant.

The culprit is signal to noise ratio.

By the end of the talk, we will see that the sensitivity is due to

- High variability of the outcome (household profit) relative to
- A small signal driving the conclusion (statistical significance)

Dropping data: Motivation

Would you be concerned if you could **reverse your conclusion** by removing a **small proportion** (say, 0.1%) of your data?

Dropping data: Motivation

Would you be concerned if you could **reverse your conclusion** by removing a **small proportion** (say, 0.1%) of your data?

Not always! But sometimes, surely yes.

Thinking without random noise can be helpful.

Suppose you have a farm, and want to know whether your average yield is greater than 170 bushels per acre. At harvest, you measure 200 bushels per acre.

- Scenario one: If your yield is greater than 170 bushels per acre, you make a profit.
 - Don't care about sensitivity to small subsets
- Scenario two: You want to recommend your farming methods to a friend across the valley.
 - Might care about sensitivity to small subsets

For example, often in economics:

- Small fractions of data are missing not-at-random,
- Policy population is different from analyzed population,
- We report a convenient summary (e.g. mean) of a complex effect,
- Models are stylized proxies of reality.

Question 1:

How do we find influential datapoints?

Which estimators do we study?

Z-estimators. Suppose we have N data points $\vec{d} = d_1, \dots, d_N$. Then:

$$\hat{\theta} := \vec{\theta} \text{ such that } \sum_{n=1}^N G(\vec{\theta}, d_n) = 0_P.$$

Examples: MLE, OLS, VB, &c (all minimizers of smooth empirical loss).

Function of interest. Qualitative decision based on $\phi(\hat{\theta}) \in \mathbb{R}$. E.g.:

- A particular component: $\phi(\theta) = \theta_d$
- The end of a confidence interval: $\phi(\theta) = \theta_d + \frac{1.96}{\sqrt{N}} \hat{\sigma}(\hat{\theta})$

Fix a proportion $0 < \alpha \ll 1$ of points to drop and find a set $\mathcal{S} \subset \{1, \dots, N\}$ with $|\mathcal{S}| \leq \lfloor \alpha N \rfloor$ that extremizes $\phi(\hat{\theta})$ when dropped.

- **Problem:** There are many sets with $|\mathcal{S}| \leq \lfloor \alpha N \rfloor$.
 - E.g., in Angelucci et al. [2015], $\binom{16,560}{15} \approx 1.5 \cdot 10^{51}$
- **Problem:** Evaluating $\phi(\hat{\theta}(\vec{d}_{-\mathcal{S}}))$ requires an estimation problem.
 - E.g., in Angelucci et al. [2015] computing the OLS estimator.
 - Other examples are even harder (VB, machine learning)

An approximation is needed!

Which estimators do we study?

Suppose we have N data points d_1, \dots, d_N . Then:

$$\hat{\theta} := \vec{\theta} \text{ such that } \sum_{n=1}^N G(\vec{\theta}, d_n) = 0_P.$$

Leave points out by setting their elements of \vec{w} to zero.

Which estimators do we study?

Suppose we have N data points d_1, \dots, d_N . Then:

$$\hat{\theta}(\vec{w}) := \vec{\theta} \text{ such that } \sum_{n=1}^N \vec{w}_n G(\vec{\theta}, d_n) = 0_P.$$

Leave points out by setting their elements of \vec{w} to zero.

Taylor series approximation.

Is there a \vec{w} , with $\lfloor \alpha N \rfloor$ zeros, such that $\phi(\hat{\theta}(\vec{w})) - \phi(\hat{\theta}) \geq \Delta$?

To simplify the search over \vec{w} , we form the Taylor series approximation:

$$\phi(\hat{\theta}(\vec{w})) - \phi(\hat{\theta}) \approx \phi^{\text{lin}}(\vec{w}) - \phi(\hat{\theta}) := - \sum_{n: \vec{w}_n=0} \psi_n, \text{ where } \psi_n := \left. \frac{\partial \phi(\hat{\theta}(\vec{w}))}{\partial \vec{w}_n} \right|_{\vec{1}}.$$

Taylor series approximation.

Is there a \vec{w} , with $\lfloor \alpha N \rfloor$ zeros, such that $\phi(\hat{\theta}(\vec{w})) - \phi(\hat{\theta}) \geq \Delta$?

To simplify the search over \vec{w} , we form the Taylor series approximation:

$$\phi(\hat{\theta}(\vec{w})) - \phi(\hat{\theta}) \approx \phi^{\text{lin}}(\vec{w}) - \phi(\hat{\theta}) := - \sum_{n: \vec{w}_n=0} \psi_n, \text{ where } \psi_n := \left. \frac{\partial \phi(\hat{\theta}(\vec{w}))}{\partial \vec{w}_n} \right|_{\vec{1}}.$$

The values ψ_n are the “**empirical influence function.**” [?]

The ψ_n can be **easily and automatically** computed from $\hat{\theta}$.

The approximation is **typically accurate** for small α .

Taylor series approximation.

Is there a \vec{w} , with $\lfloor \alpha N \rfloor$ zeros, such that $\phi(\hat{\theta}(\vec{w})) - \phi(\hat{\theta}) \geq \Delta$?

To simplify the search over \vec{w} , we form the Taylor series approximation:

$$\phi(\hat{\theta}(\vec{w})) - \phi(\hat{\theta}) \approx \phi^{\text{lin}}(\vec{w}) - \phi(\hat{\theta}) := - \sum_{n: \vec{w}_n=0} \psi_n, \text{ where } \psi_n := \left. \frac{\partial \phi(\hat{\theta}(\vec{w}))}{\partial \vec{w}_n} \right|_{\vec{1}}.$$

The values ψ_n are the “**empirical influence function.**” [?]

The ψ_n can be **easily and automatically** computed from $\hat{\theta}$.

The approximation is **typically accurate** for small α .

Is there a \vec{w} , with $\lfloor \alpha N \rfloor$ zeros, such that $\phi^{\text{lin}}(\vec{w}) - \phi(\hat{\theta}) \geq \Delta$?

Taylor series approximation.

Is there a \vec{w} , with $\lfloor \alpha N \rfloor$ zeros, such that $\phi(\hat{\theta}(\vec{w})) - \phi(\hat{\theta}) \geq \Delta$?

To simplify the search over \vec{w} , we form the Taylor series approximation:

$$\phi(\hat{\theta}(\vec{w})) - \phi(\hat{\theta}) \approx \phi^{\text{lin}}(\vec{w}) - \phi(\hat{\theta}) := - \sum_{n: \vec{w}_n=0} \psi_n, \text{ where } \psi_n := \left. \frac{\partial \phi(\hat{\theta}(\vec{w}))}{\partial \vec{w}_n} \right|_{\vec{1}}.$$

The values ψ_n are the “**empirical influence function.**” [?]

The ψ_n can be **easily and automatically** computed from $\hat{\theta}$.

The approximation is **typically accurate** for small α .

Is there a \vec{w} , with $\lfloor \alpha N \rfloor$ zeros, such that $\phi^{\text{lin}}(\vec{w}) - \phi(\hat{\theta}) \geq \Delta$?

Easy! The most influential points for $\phi^{\text{lin}}(\vec{w})$ have the most negative ψ_n .

Taylor series approximation.

Procedure:

Taylor series approximation.

Procedure:

- 1 Compute the “original” estimator, $\hat{\theta}$ and $\phi(\hat{\theta})$.

Procedure:

- 1 Compute the “original” estimator, $\hat{\theta}$ and $\phi(\hat{\theta})$.
- 2 Compute and sort the influence scores, $\psi_{(1)} \leq \psi_{(2)} \leq \dots \leq \psi_{(N)}$.

Procedure:

- 1 Compute the “original” estimator, $\hat{\theta}$ and $\phi(\hat{\theta})$.
- 2 Compute and sort the influence scores, $\psi_{(1)} \leq \psi_{(2)} \leq \dots \leq \psi_{(N)}$.
- 3 Let \vec{w}^* leave out the data corresponding to $\psi_{(1)}, \dots, \psi_{(\lfloor \alpha N \rfloor)}$.

Procedure:

- 1 Compute the “original” estimator, $\hat{\theta}$ and $\phi(\hat{\theta})$.
- 2 Compute and sort the influence scores, $\psi_{(1)} \leq \psi_{(2)} \leq \dots \leq \psi_{(N)}$.
- 3 Let \vec{w}^* leave out the data corresponding to $\psi_{(1)}, \dots, \psi_{(\lfloor \alpha N \rfloor)}$.
- 4 Report non-robustness if $\Delta \leq \phi^{\text{lin}}(\vec{w}^*) - \phi(\hat{\theta}) = -\sum_{n=1}^{\lfloor \alpha N \rfloor} \psi_{(n)}$.

Procedure:

- 1 Compute the “original” estimator, $\hat{\theta}$ and $\phi(\hat{\theta})$.
- 2 Compute and sort the influence scores, $\psi_{(1)} \leq \psi_{(2)} \leq \dots \leq \psi_{(N)}$.
- 3 Let \vec{w}^* leave out the data corresponding to $\psi_{(1)}, \dots, \psi_{(\lfloor \alpha N \rfloor)}$.
- 4 Report non-robustness if $\Delta \leq \phi^{\text{lin}}(\vec{w}^*) - \phi(\hat{\theta}) = -\sum_{n=1}^{\lfloor \alpha N \rfloor} \psi_{(n)}$.
- 5 **Optional:** Compute $\hat{\theta}(\vec{w}^*)$, and verify that $\Delta \leq \phi(\hat{\theta}(\vec{w}^*)) - \phi(\hat{\theta})$.

Computing the influence function.

How to compute $\psi_n := \left. \frac{\partial \phi(\hat{\theta}(\vec{w}))}{\partial \vec{w}_n} \right|_{\vec{1}}$? Recall $\sum_{n=1}^N \vec{w}_n G(\hat{\theta}(\vec{w}), d_n) = 0_P$.

Step zero: Implement software to compute $G(\theta, d_n)$ and $\phi(\theta)$. Find $\hat{\theta}$.

Step one: By the chain rule, $\psi_n = \left. \frac{\partial \phi(\hat{\theta}(\vec{w}))}{\partial \vec{w}_n} \right|_{\vec{1}} = \left. \frac{d\phi(\theta)}{d\theta^T} \right|_{\hat{\theta}} \left. \frac{\partial \hat{\theta}(\vec{w})}{\partial \vec{w}_n} \right|_{\vec{1}}$.

Step two: By the implicit function theorem:

$$\left. \frac{\partial \hat{\theta}(\vec{w})}{\partial \vec{w}_n} \right|_{\vec{1}} = \frac{1}{N} \left(\frac{1}{N} \sum_{n'=1}^N \left. \frac{\partial}{\partial \theta^T} G(\vec{\theta}, d_{n'}) \right|_{\hat{\theta}} \right)^{-1} G(\hat{\theta}, d_n).$$

Step three: Use *automatic differentiation* on $\phi(\theta)$ and $G(\theta, d_n)$ from step zero to compute $\left. \frac{\partial \phi(\theta)}{\partial \theta^T} \right|_{\hat{\theta}}$ and $\left. \frac{\partial}{\partial \theta^T} G(\vec{\theta}, d_n) \right|_{\hat{\theta}}$.

-
- The user does step zero. The rest is automatic.
 - The primary computational expense is the Hessian inverse.
 - Automatic differentiation is the chain rule applied to a program.
 - Typically $\psi_n = O(N^{-1})$.

Question 2:

What makes an estimator non-robust?

Question 3:

When is our approximation accurate?

Conclusion: Related work and future directions

Tamara Broderick, Ryan Giordano, Rachael Meager (alphabetical authors)
“An Automatic Finite-Sample Robustness Metric: Can Dropping a Little Data Change Conclusions?”

<https://arxiv.org/abs/2011.14999>

M. Angelucci, D. Karlan, and J. Zinman. Microcredit impacts: Evidence from a randomized microcredit program placement experiment by Compartamos Banco. *American Economic Journal: Applied Economics*, 7(1):151–82, 2015.